

Agenda

- Motivations: what question(s) can info theory help us answer?
- Relationship diagram between key quantities
- Problems

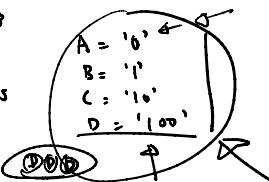
① Find min # bits on avg to send msg to someone

$H(X)$

↓ 4 letter language : A, B, C, D
 '0' '1' '10' '11' ≈ 2.6 bits

What if we knew, AAC AAB D AAA

$$\begin{aligned}A &= 0.8 \\B &= 0.1 \\C &= 0.05 \\D &= 0.05\end{aligned}$$



$$\begin{aligned}A &= '0' \\B &= '1' \\C &= '10' \\D &= '100'\end{aligned}$$

0100

ABA or AD

Source Coding Theorem

a bunch of symbols, e.g. A, B, C, D, E, ...
 mapping \rightarrow bit sequences

$X_1, X_2, \dots, X_n \sim i.i.d P_X$, $l(x) :=$ length of x (in bits)

\exists a source coding scheme, s.t.

$$\lim_{n \rightarrow \infty} \mathbb{E} \left[\frac{1}{n} l(X_1, \dots, X_n) \right] \leq H(X) + \epsilon$$

$H(X)$

$$\forall \epsilon > 0, P(|X_n - X| > \epsilon) \leq \frac{\text{Var}(X)}{\epsilon}$$

$$H(X) = \mathbb{E} \left[\log \frac{1}{P_X(x)} \right] = \sum_x P_X(x) \log \frac{1}{P_X(x)}$$

$$P_X(x) = 0.002$$

$P_X(x)$ = really big

$$\log \frac{1}{P_X(x)}$$

Expected 'surprise'

'uncertainty'

'info content'

$$H(X), H(X,Y), H(X|Y), I(X;Y)$$

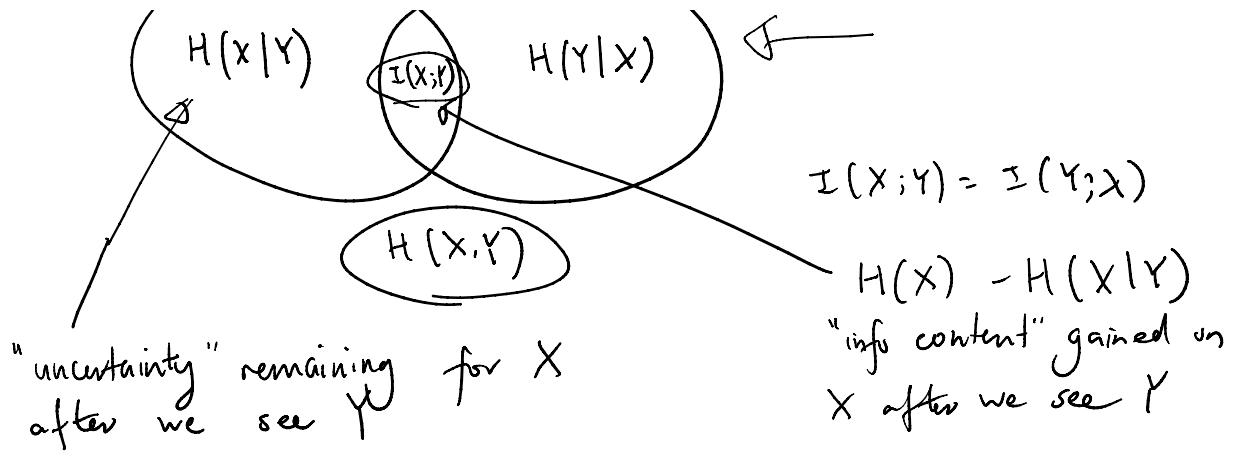
$H(X)$

$H(Y)$

$H(X|Y)$

$H(Y|X)$

$I(X;Y)$



1. Mutual Information and Noisy Typewriter

The **mutual information** of X and Y is defined as

$$I(X;Y) := H(X) - H(X|Y)$$

Here, $H(X|Y)$ denotes the **conditional entropy** of X given Y , which is defined as:

$$\begin{aligned} \underline{H(X|Y)} &= \sum_{y \in \mathcal{Y}} p_Y(y) H(X|Y=y) \\ &= \sum_{y \in \mathcal{Y}} p_Y(y) \sum_{x \in \mathcal{X}} p_{X|Y}(x|y) \log_2 \frac{1}{p_{X|Y}(x|y)} \\ &= - \sum_{y \in \mathcal{Y}, x \in \mathcal{X}} p_{X,Y}(x,y) \log_2 \frac{p_{X,Y}(x,y)}{p_Y(y)} \end{aligned}$$

The interpretation of conditional entropy is the average amount of uncertainty remaining in the random variable X after observing Y . The interpretation of mutual information is therefore the amount of information about X gained by observing Y .

- (a) Show that $\underline{H(X,Y)} = H(Y) + \underline{H(X|Y)} = H(X) + H(Y|X)$. This is often called the **Chain Rule**. Interpret this rule.

$$\begin{aligned} H(X|Y) &= \mathbb{E} \left[\log \frac{1}{p_{X|Y}(X|Y)} \right] \\ \underline{H(X,Y)} &= \mathbb{E} \left(\log \frac{1}{p_{X,Y}(X,Y)} \right) \\ &= \mathbb{E} \left[\log \frac{1}{p_{X|Y}(X|Y)} \right] + \left[\log \frac{1}{p_Y(Y)} \right] \\ &= \mathbb{E} \left(\log \frac{1}{p_{X|Y}(X|Y)} \right) + \mathbb{E} \left(\frac{1}{p_Y(Y)} \right) \\ &= H[X|Y] + H(Y) \end{aligned}$$

- (b) Show that $\underline{I(X;Y)} = H(X) + H(Y) - H(X,Y)$. Note that this shows that $I(X;Y) = I(Y;X)$, i.e., mutual information is symmetric.

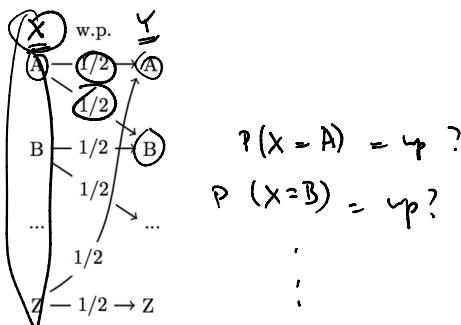
$$\begin{aligned} \underline{I(X;Y)} &= H(X) - \underline{H(X|Y)} = H(X) - (H(X,Y) - H(Y)) \\ &\approx H(X) + H(Y) - H(X,Y) \end{aligned}$$

$$\begin{aligned} &= H(X) + H(Y) - H(X, Y) \\ &= H(Y) - H(Y|X) \end{aligned}$$

$$I(X;Y) = H(X) - \underline{H(X|Y)}$$

$$\cancel{H(X,Y)} = \cancel{H(X|Y)}$$

(c) Consider the noisy typewriter.



Each symbol gets sent to one of the adjacent symbols with probability 1/2. Let X be the input to the noisy typewriter, and let Y be the output (X is a random variable that takes values in the English alphabet). What is the distribution of X that maximizes $I(X;Y)$?

Note: distribution that maximizes the entropy of a random variable is the uniform distribution

$$\begin{aligned} I(X;Y) &= (H(X) - H(X|Y)) \\ &= H(Y) - H(Y|X) = 1 \end{aligned}$$

ans. Let X be uniform on $\{A, B, \dots, Z\}$

$$Z \sim \text{Bern}(\frac{1}{2})$$

$$\begin{cases} X = A \quad \text{w.p. } \frac{1}{26} \\ \vdots \\ X = Z \quad \text{w.p. } \frac{1}{26} \end{cases} \Rightarrow Y = \begin{cases} A \quad \text{w.p. } \frac{1}{26} \\ \vdots \\ Z \quad \text{w.p. } \frac{1}{26} \end{cases}$$

Note It turns out that $I(X;Y) \geq 0$ with equality if and only if X and Y are independent. The mutual information is an important quantity for channel coding.

$$\begin{aligned} H(Z) &= \sum_z p_z(z) \log \frac{1}{p_z(z)} \\ &= \frac{1}{2} \log_2 2 + \frac{1}{2} \log_2 2 \\ &= \frac{1}{2} + \frac{1}{2} = 1 \end{aligned}$$

$$Y|X \sim \text{Bern}(\frac{1}{2})$$

$\Rightarrow I(X;Y)$ maximised

2. Entropy of a Sum
- (a) Let X_1, X_2 be i.i.d. Bernoulli(1/2) (fair coin flips). Calculate $H(X_1 + X_2)$ and show that $H(X_1 + X_2) \geq \cancel{H(X_1)}$

$$X_1 + X_2 \sim \text{Bin}(2, \frac{1}{2})$$

$$Z = X_1 + X_2 = \begin{cases} 0 & \text{w.p. } \frac{1}{4} \\ 1 & \text{w.p. } \frac{1}{2} \\ 2 & \text{w.p. } \frac{1}{4} \end{cases}$$

$$x_1 + x_2 \sim \text{Bin}(2, \frac{1}{2})$$

$$Z = X_1 + X_2 = \begin{cases} 1 & \text{if } \frac{1}{2} \\ 2 & \text{if } \frac{1}{4} \end{cases}$$

$$\begin{aligned} H(X_1 + X_2) &= \sum_z P_{X_1+X_2}(z) \log \frac{1}{P_{X_1+X_2}(z)} \\ &= \frac{1}{4} \log_2 4 + \frac{1}{2} \log_2 2 + \frac{1}{4} \log_2 4 \\ &\approx \frac{1}{2} + \frac{1}{2} + \frac{1}{2} = \frac{3}{2} \end{aligned}$$

$$H(X_1) = I \quad (\text{by formula})$$

- (b) It turns out that in general adding independent random variables increases entropy. To prove this, first let's prove that conditioning decreases entropy. Let X and Y be two (not necessarily i.i.d) random variables, prove that $H(X|Y) \leq H(X)$.
Hint: Use the fact that $I(X; Y) \geq 0$ for any r.v.s X and Y .

$$I(X; Y) = \underbrace{H(X)}_{\text{H}} - \underbrace{H(X|Y)}_{\text{H}} \geq 0$$

- (c) Now suppose X_1, X_2 are two independent (but not necessarily identically distributed) random variables. Prove that $H(X_1 + X_2) \geq H(X_1)$.

$$\begin{aligned} H(X_1 + X_2) &\geq H(X_1 + X_2 | X_2) = H(X_1 | X_2) = H(X_1) \\ &= \sum_{z, x_2} P_{X_1+X_2, X_2}(z, x_2) \log \frac{1}{P_{X_1+X_2 | X_2}(z | x_2)} \end{aligned}$$